

LAW OFFICES

AMIN, TUROCY, AND CALVIN LLP

24TH FLOOR, NATIONAL CITY CENTER

1900 EAST NINTH STREET

CLEVELAND, OHIO 44114

TELEPHONE: 216-696-8730

FACSIMILE: 216-696-8731

EMAIL: BRAYAPROLU@THEPATENTATTORNEYS.COM

THIS MESSAGE IS INTENDED ONLY FOR THE USE OF THE INDIVIDUAL OR ENTITY TO WHICH IT IS ADDRESSED AND MAY CONTAIN INFORMATION THAT IS PRIVILEGED, CONFIDENTIAL AND EXEMPT FROM DISCLOSURE UNDER APPLICABLE LAW. IF THE READER OF THIS MESSAGE IS NOT THE INTENDED RECIPIENT OR THE EMPLOYEE OR AGENT RESPONSIBLE FOR DELIVERING THE MESSAGE TO THE INTENDED RECIPIENT, YOU ARE HEREBY NOTIFIED THAT ANY DISSEMINATION, DISTRIBUTION OR COPYING OF THIS COMMUNICATION IS STRICTLY PROHIBITED. IF YOU HAVE RECEIVED THIS COMMUNICATION IN ERROR, PLEASE NOTIFY US IMMEDIATELY, AND RETURN THE ORIGINAL MESSAGE TO US AT THE ADDRESS LISTED BELOW VIA UNITED STATES MAIL. THANK YOU.

Date: May 28, 2008

TO: Scott M. Sciacca – United States Patent and Trademark Office

FROM: Bhavani S. Rayaprolu

In re patent application of:

Applicants: Andrew Laucius, *et al.*

Serial No: 10/750,011

Filing Date: December 31, 2003

Examiner: Scott M. Sciacca

Art Unit: 2146

Title: INCREMENTAL WEB CRAWLER USING CHUNKS

TOTAL NUMBER OF PAGES (INCLUDING THIS PAGE): 6

PROPOSED AMENDMENTS

Listing of Claims:

1. (Previously Presented) A system that facilitates incremental web crawls comprising:
an indexer that places items with similar properties into respective chunks; and,
a chunk map that stores at least some of the properties associated with the respective chunk, wherein the properties are at least one of average time between change or average importance of documents comprising a particular chunk, the chunk map employed to facilitate an incremental web re-crawl.
2. (Original) The system of claim 1, the items comprising information associated with a Uniform Resource Locator.
3. (Original) The system of claim 1, the items comprising at least one of an HTML file, a PDF file, a PS file, a PPT file, an XLS file and a DOC file.
4. (Original) The system of claim 1, the items receives from a crawler, the crawler responsible for a specific set of Uniform Resource Locators.
5. (Original) The system of claim 1, further comprising a master control process that can modify the chunk map to facilitate load balancing amongst a plurality of crawlers.
6. (Original) The system of claim 1, further comprising a master control process that serves as an interface between a crawler and a re-crawl controller.
7. (Original) The system of claim 6, wherein the master control process maintains a known chunks table that stores information for components of a system.

8. (Original) The system of claim 6, wherein the master control process exposes an interface for communication with a component of the system.
9. (Original) The system of claim 8, wherein the interface returns a list of chunks the component should have and where to get the chunks.
10. (Original) The system of claim 8, wherein the interface returns a list of the chunks that should be actively served by the component.
11. (Original) The system of claim 8, wherein the interface returns a range of chunk identifiers to use in building a new chunk by the component.
12. (Original) The system of claim 8, wherein the interface causes an old chunk to be retired by the system.
13. (Original) The system of claim 6, wherein the master control process facilitates movement of chunks from one component to another component.
14. (Original) The system of claim 13, wherein movement of chunks is based, at least in part, upon at least one of rebalancing index servers after one goes down, re-crawling pages previously crawled, and, restoring a state of a crawler after it has crashed.
15. (Original) The system of claim 1, further comprising a re-crawl component that employs the chunk map to determine which chunks, if any, to re-crawl at a particular time.
16. (Cancelled)
17. (Original) The system of claim 1, further comprising an index chunk that stores information associated with an index of at least some of the items.

18. (Original) The system of claim 1, further comprising a rank chunk that stores a static rank associated with an index chunk.
19. (Currently Amended) A method of performing document re-crawl comprising:
parsing a first chunk for uniform resource locators, wherein ~~the a chunk map[[s]] that stores properties associated with the respective chunk stored in a chunk table are is~~ employed to determine the first chunk;
re-crawling the uniform resource locators; and,
forming a second chunk based, at least in part, upon the re-crawled uniform resource locators.
20. (Original) The method of claim 19 comprising at least one of the following acts:
determining whether any chunks are to be retired;
moving the first chunk; and,
destroying the first chunk.
21. (Original) One or more computer readable media having stored thereon computer executable instructions for carrying out the method of claim 19.
22. (Previously Presented) A method of performing document re-crawl comprising:
accessing a chunk map containing properties associated with respective chunks of data as a result of one or more web crawls, wherein the properties are at least one of average time between change or average importance of documents comprising a particular chunk; and,
periodically determining, based on the properties in the chunk map, whether to re-crawl one or more of the chunks of data.
23. (Original) The method of claim 22, the period determination being based, at least in part, upon, at least one of average time between change and average importance of documents comprising a particular chunk.

24. (Currently Amended) A data packet transmitted between two or more computer components that facilitates document re-crawl, the data packet comprising:
 - a chunk header that includes metadata associated with the data packet, wherein a chunk comprises document files that have similar properties;
 - an offset section that provides offset information associated with document files; and,
 - the document files that include content found on the Internet, wherein the average of the at least one of the properties of all the document files determines if the document should be re-crawled.
25. (Original) The data packet of claim 24, at least one of the document files comprising at least one of an HTML file, a PDF file, a PS file, a PPT file, an XLS file and a DOC file.
26. (Previously Presented) A system that facilitates increment web crawls comprising:
 - means for placing items with similar properties into respective chunks; and,
 - means for storing at least some of the properties associated with the respective chunk, wherein the properties are at least one of average time between change or average importance of documents comprising a particular chunk, and employing the properties to facilitate an incremental web re-crawl.
27. (Original) The system of claim 26, the items comprising information associated with a Uniform Resource Locator.
28. (Original) The system of claim 26, the items comprising at least one of an HTML file, a PDF file, a PS file, a PPT file, an XLS file and a DOC file.

REMARKS

At the cited portions, Evans discloses storing auxiliary information pertaining to the encountered web sites in a database. This stored information is employed to determine how often to recrawl. Thus, the system stores the properties of each and every document in a database. In contrast, the claimed invention allows for chunking, wherein a set of documents that have similar properties are placed in a chunk. Each chunk and its properties are stored in a chunk map and this allows for the chunk to be manipulated as one set. Each document does not have to be analyzed, but the whole chunk can be recrawled depending on its properties. Thus, Evans is silent regarding *a chunk map that stores at least some of the properties associated with the respective chunk, and the chunk map employed to facilitate an incremental web re-crawl* as recited by independent claims 1 and 26.

Independent claim 19 recites *a method of performing document re-crawl comprising: parsing a first chunk for uniform resource locators, wherein a chunk map that stores properties associated with the respective chunk is employed to determine the first chunk; re-crawling the uniform resource locators; and forming a second chunk based, at least in part, upon the re-crawled uniform resource locators.* Najork *et al.* and Evans are silent regarding such novel features. At the cited portions, Najork *et al.* discloses performing a recrawl of a queue and enqueueing any new URL's into the front-end queue or in another queue depending on host identifier of the new URL. In contrast, the claimed invention allows for *forming a second chunk*, if the indexer determines that the document belonging to the new URL does not belong to an existing chunk. Thus Najork *et al.* is silent regarding *forming a second chunk based, at least in part, upon the re-crawled uniform resource locators* as recited by independent claim 19. Evans does not compensate for the aforementioned deficiency of Najork *et al.*